

```
INSERT INTO RiIDF(token, idf)
SELECT  T.token, LOG(S.size)-LOG(COUNT(UNIQUE(*)))
FROM    RiTokens T, RiSize S
GROUP BY T.token, S.size
```

(a) Relation with token *idf* counts

```
INSERT INTO RiLength(tid, len)
SELECT  T.tid, SQRT(SUM(I.idf*I.idf*T.tf*T.tf))
FROM    RiIDF I, RiTF T
WHERE   I.token = T.token
GROUP BY T.tid
```

(c) Relation with weight-vector lengths

```
INSERT INTO RiSum(token, total)
SELECT  R.token, SUM(R.weight)
FROM    RiWeights R
GROUP BY R.token
```

(e) Relation with total token weights

```
INSERT INTO RiTF(tid, token, tf)
SELECT  T.tid, T.token, COUNT(*)
FROM    RiTokens T
GROUP BY T.tid, T.token
```

(b) Relation with token *tf* counts

```
INSERT INTO RiWeights(tid, token, weight)
SELECT  T.tid, T.token, I.idf*T.tf/L.len
FROM    RiIDF I, RiTF T, RiLength L
WHERE   I.token = T.token AND T.tid = L.tid
```

(d) Final relation with normalized tuple weight vectors

```
INSERT INTO RiSize(size)
SELECT  COUNT(*)
FROM    Ri
```

(f) Dummy relation used to create *RiIDF*

Fig. 1

```
SELECT r1w.tid AS tid1, r2w.tid AS tid2
FROM R1Weights r1w, R2Weights r2w
WHERE r1w.token = r2w.token
GROUP BY r1w.tid, r2w.tid
HAVING SUM(r1w.weight*r2w.weight) ≥  $\phi$ 
```

Fig. 2

```
SELECT rw.tid, rw.token, rw.weight/rs.total AS P
FROM   RiWeights rw, RiSum rs
WHERE  rw.token = rs.token
```

Fig 3

```
INSERT INTO RiSample(tid,token,c)
SELECT  rw.tid, rw.token, ROUND(S * rw.weight/rs.total, 0) AS c
FROM    RiWeights rw, RiSum rs
WHERE   rw.token = rs.token
```

Fig. 4

```
SELECT  riw.tid AS tid1, r2s.tid AS tid2
FROM    R1weights riw, R2sample r2s, R2sum r2sum, R1V r1v
WHERE   riw.token = r2s.token AND riw.token = r2sum.token AND riw.tid = r1v.tid
GROUP BY riw.tid, r2s.tid, r1v.Tv
HAVING  SUM(riw.weight * r2sum.total / r1v.Tv)  $\geq S * \phi' / r1v.Tv$ 
```

Fig. 5

```
SELECT tid1, tid2
FROM
(
  SELECT  r1w.tid AS tid1, r2s.tid AS tid2, SUM(r1w.weight * r2sum.total) AS Ci
  FROM    R1weights r1w, R2sample r2s, R2sum r2sum
  WHERE   r1w.token = r2s.token AND r1w.token = r2sum.token AND r1w.tid = r1v.tid
  GROUP BY r1w.tid, r2s.tid
  UNION ALL
  SELECT  r1s.tid AS tid1, r2w.tid AS tid2, SUM(r2w.weight * r1sum.total) AS Ci
  FROM    R2weights r2w, R1sample r1s, R1sum r1sum
  WHERE   r2w.token = r1s.token AND r2w.token = r1sum.token AND r2w.tid = r2v.tid
  GROUP BY r2w.tid, r1s.tid
) SYM
GROUP BY tid1, tid2
HAVING AVG(Ci)  $\geq S * \phi'$ 
```

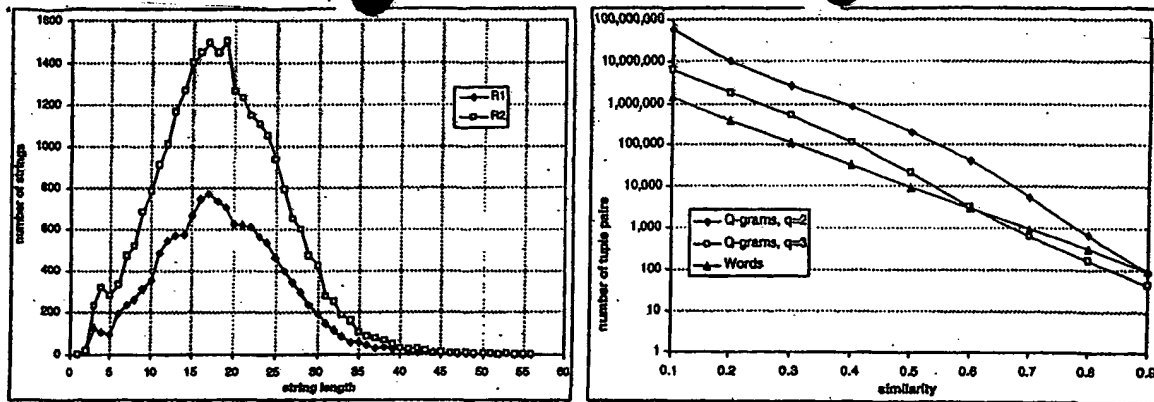
Fig. 6

7/12

```
SELECT  r1s.tid AS tid1, r2s.tid AS tid2
FROM    R1Sample r1s, R2Sample r2s, R1Sum r1sum, R2Sum r2sum
WHERE   r1s.token = r1sum.token AND R2Sample.token = r2sum.token AND r1s.token = r2s.token
GROUP BY r1s.tid, r2s.tid
HAVING  SUM(r1sum.total * r2sum.total)  $\geq S * S * \phi'$ 
```

Fig. 7

8/12



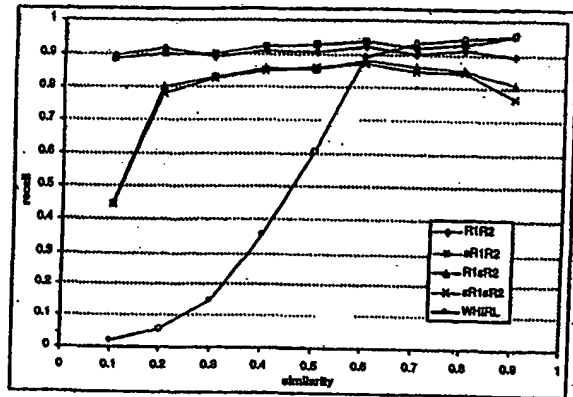
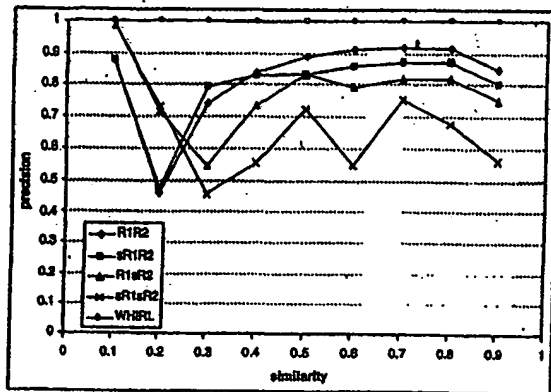
(a) String lengths in data sets  $R_1$  and  $R_2$ .

(b) The size of  $R_1 \bowtie_{\theta} R_2$  for different similarity thresholds and token choices.

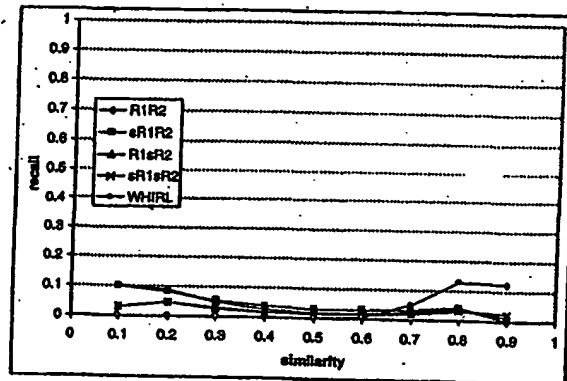
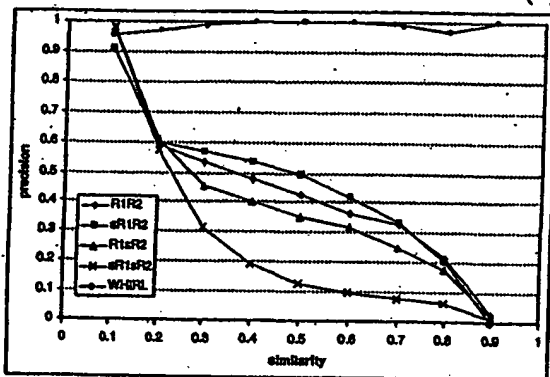
Fig 8



9/12



(a) Words



(b) Q-grams with  $q=2$

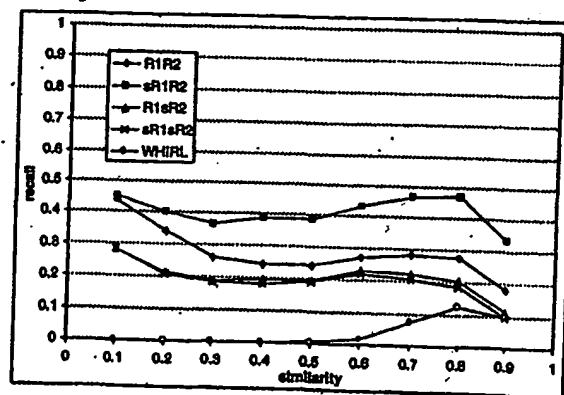
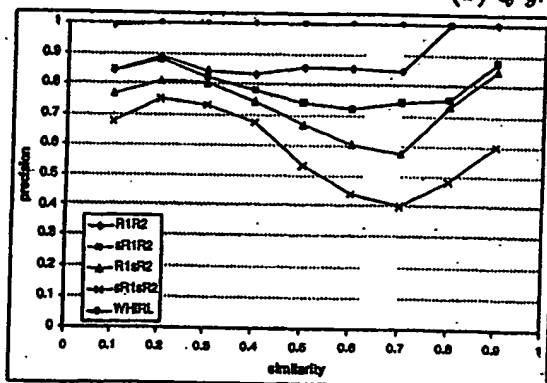


Fig. 9

10/12

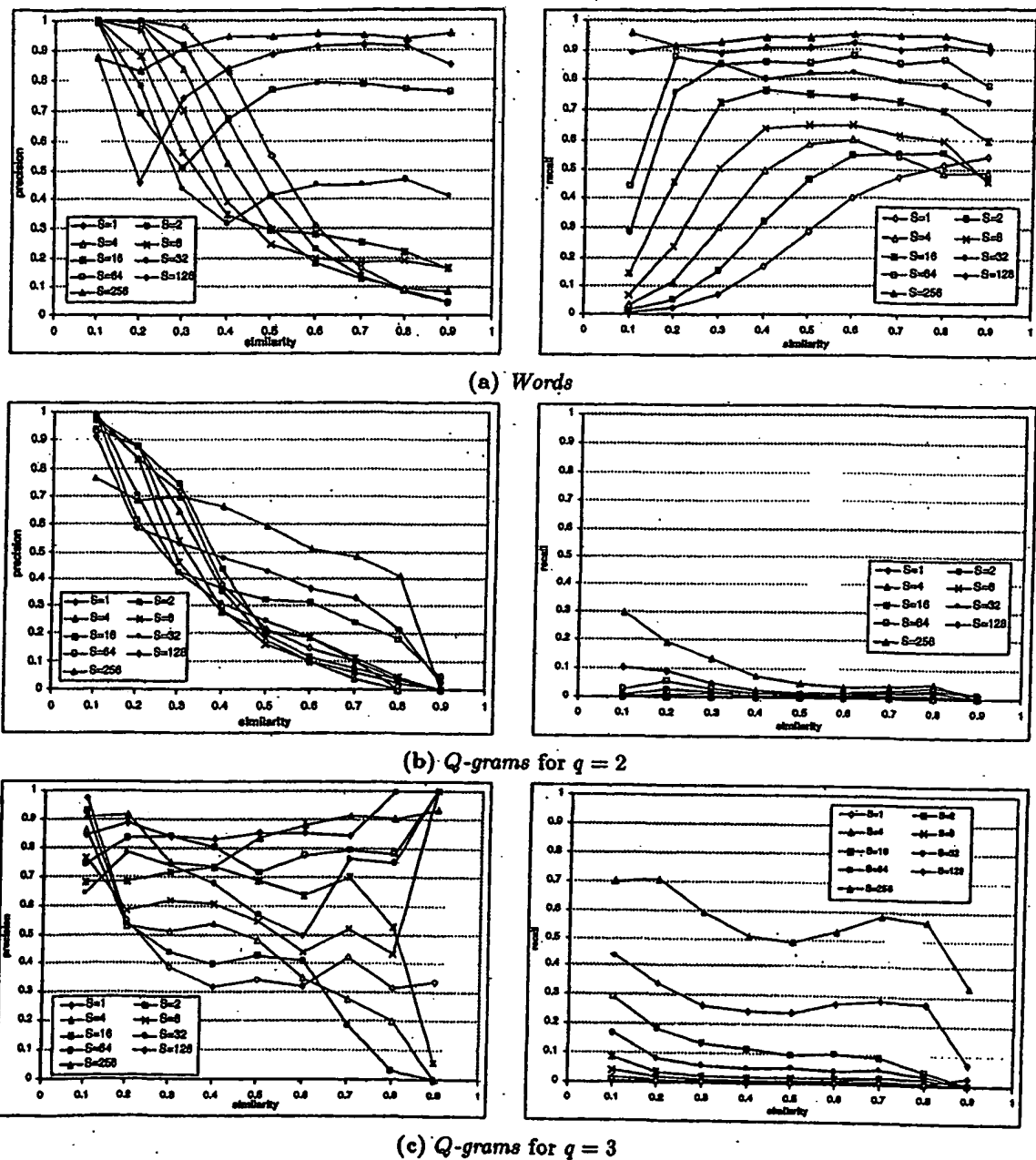


Fig. 10

11/12

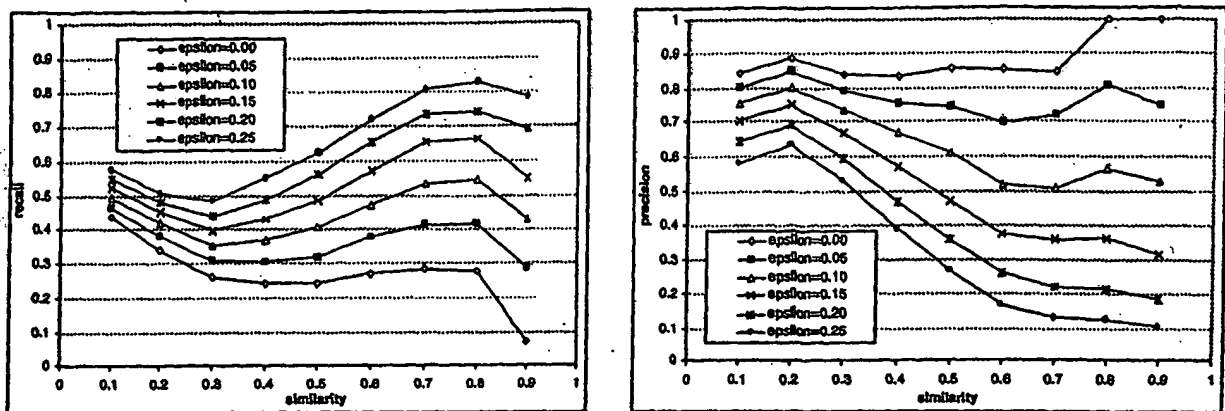
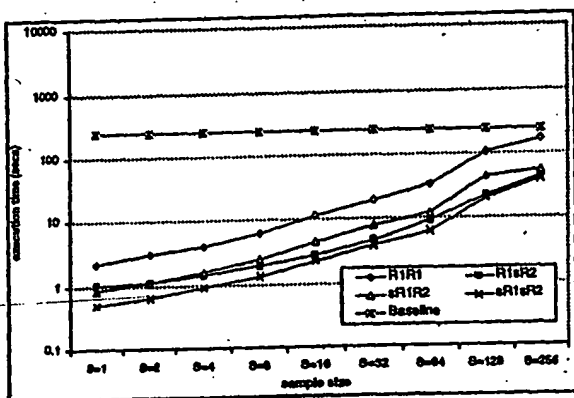
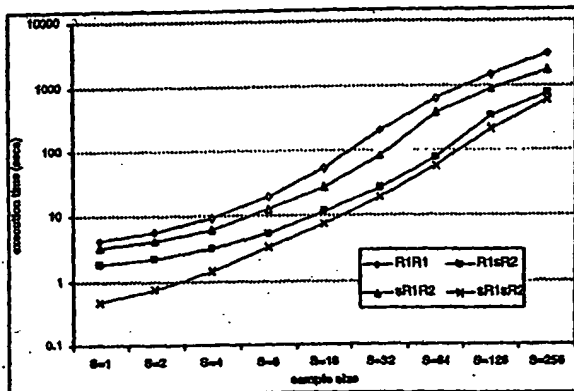


Fig. 11

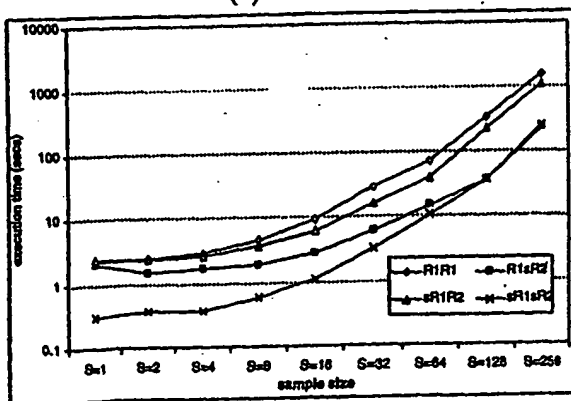
12/12



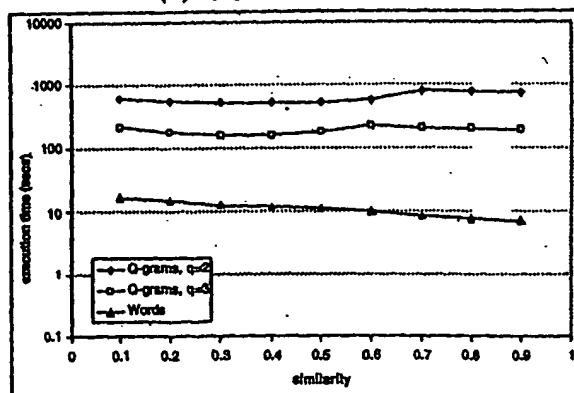
(a) Words



(b) Q-grams with  $q=3$



(c) Q-grams with  $q=2$



(d) WHIRL

Fig. 12